

K-Edge Component의 그래픽 정의를 이용한 집합화 클러스터링

임 근
서울보건대학 전산정보처리과
e-mail:LK04@www.shjc.ac.kr

Aggregation Clustering using Graphic Conecpt of K- Edge Component

Lim, Keun
Dept of Computer Information Process, Seoul Health Junior
College

요약

본 논문에서는 클러스터 정의시 사용하게 되는 특성으로 노드간 패스 수에 기반한 K-edge 컴포넌트의 그래픽 정의 방법과 노드를 클러스터화 하는 집합화(Aggregation) 방법을 제시하였다. 집합화된 하이퍼텍스트 분리를 통해 이전 결과를 개선할 수 있으며, 집합내의 노드 간 관련성을 가시화하여 비교할 수 있다.

1. 서론

클러스터링 알고리즘은 사실과 현상의 집합에서 여러 구성요소의 구조를 찾아보기 위해서 사용되었다. 특히 자연과학분야에서 분류와 조직을 통한 데이터의 검색 분야에 적용하는 예가 많이 있다. 본 논문에서는 Term 클러스터링과 구문 문법형태를 통한 자연어 문장을 번역하기 위한 방법을 제시하고, 또한 클러스터 정의에 사용하게 되는 특성으로 노드 간 패스 수에 기반한 K-edge 컴포넌트의 그래픽 정의를 이용하여, 노드를 클러스터화 하는 집합화(Aggregation) 방법을 제시한다. 2가지 방법은 모두 도큐먼트 검색에 대한 텍스트 표현을 통해서 검색의 질적인 면을 향상시켜보려는 목적을 가지고 있다. 또한 집단화 방법을 적용함으로써 클래스들의 클러스터화를 실현해 보고자 한다.

본 논문의 구성은 2장에서 관련 연구로 클러스터링 개념과 조건 등을 설명하며, 3장에서 클러스터링 메소드, 4장에서 적용 예와 결론을 맺는다.

2. 관련 연구

2.1 Term 클러스터링 개념

Term 클러스터링은 기존에 텍스트 표현으로부터 추출한 Term을 그룹화하기 위한 클러스터 분석 중에 하나의 응용분야로 특성추출 방법이라고 할 수 있으며 패턴분류에 유용하다. Term 클러스터는 비교적 적은 빈도로 유사한 Term에 의해 제한을 받게 되며 Term 대치의 의미보다는 보충의 개념으로 이용된다. 클러스터링은 관련성 있는 도큐먼트의 분리 여부에 따라 효과가 결정된다.

2.2 클러스터링 알고리즘

클러스터링 알고리즘은 독립 개체의 자연적인 클러스터를 찾기 위한 시도로서 개체의 집합은 내부적 응집력과 외부적 결합성을 통해서 특성을 설명할 수 있다.

- Clustering Criterion : N개의 표본, 즉 $X_1, X_2, X_3, \dots, X_N$ 으로 분류한다고 가정하자. 이러한 벡터들은 임의 벡터로 표시하지 않는다. 그 이유는 벡터들이 고정적인 값이고 일반화되어 있는 요소라고 가정하기 때문이다. 각 표본은 L개 클래스

(W_1, W_2, \dots, W_N)로 표기한다. 편의상 K_i 의 값은 1과 L 사이의 정수값이라고 가정한다. 하나의 분류 요소(Classification) Ω 는 벡터로서 W_{ki} 라고 하며 Configuration X^* 는 X_i 로 나타내며 다음과 같이 정리한다.

$$\Omega = [W_{k1}, W_{k2}, \dots, W_{kn}]^T$$

$$X^* = [X_1^T, X_2^T, \dots, X_N^T]^T$$

Clustering Criterion J는 Ω 와 X^* 의 함수로 다음과 같이 표시할 수 있다.

$$J = J(W_{k1}, W_{k2}, \dots, W_{kn} : X_1^T, X_2^T, \dots, X_N^T) \\ = J(\Omega_n : X^*)$$

정의에 의하여 Classification Ω 는 다음을 만족한다.

$$J(\Omega_n : X^*) = \max_{\Omega} \text{ 또는 } \min_{\Omega} J(\Omega : X^*)$$

- Clustering algorithm : 주어진 클러스터링 문제에 대하여 Configuration X^* 는 고정적인 값이고 Classification Ω 는 변수이다. 일반적인 클러스터링 기법을 통한 방법은 실제 검색에 적용하기에는 부적합하다. 그 이유는 Classification Ω 가 이산적이고, 비순서적이기 때문이다. 그러나 Classification Ω 에서 분산에 관한 J상의 반복적 탐색 알고리즘 정의는 가능하다. 1번 반복한 Classification은 Classification $\Omega(1)$ 이라고 가정하자.

$$\Omega(1) = [W_{k1(1)}, W_{k2(1)}, \dots, W_{kn(1)}]^T$$

만일 1개의 표본이 현재 클래스 $K_{i(1)}$ 에서 클래스 j로 Clustering Criterion은 전체합 $\Delta J(i, j, 1)$ 에 의해 변환된다. 이때 $\Delta J(i, j, 1) = J(W_{k1(1)}, W_{k2(1)}, \dots, W_{kn(1)} : X^*) - J(\Omega(1) : X^*)$ 만일 $\Delta J(i, j, 1)$ 이 음수값이면, 클래스 j로의 1개 표본에 대한 재분류는 하나의 Classification을 만들게 되고, 이것은 다음 알고리즘에 기반한다.

- ① 초기 Classification $\Omega(0)$ 을 선택한다.
- ② 1번째 Classification $\Omega(1)$ 이 주어지고, $\Delta J(i, j, 1)$ 은 계산을 통해 구한다. ($j=1, 2, \dots, L ; i=1, 2, \dots, N$)
- ③ $i=1, 2, \dots, N$ 에 대하여, 1개의 표본을 클래스로 재분류한다. $\Delta J(i, j, 1) = \min_j \Delta J(i, j, 1)$ 현재 Classification에 의해 $j = k_{i(1)}$ 를 포함하고, 이 단계는 $\Omega(i+1)$ 의 형태를 취한다.
- ④ 만일 $\Omega(i+1) < \Omega(1)$ 이면, ② 단계로 간다.

이 알고리즘은 Clustering Criterion에 기반한 Classification rule의 단순한 반복 응용이라고 할 수 있다.

3. K-Component 기반 집합 클러스터링 방법
바람직한 클러스터링은 자연스럽고, 세부사항에 독

립적이며, 유사한 객체를 동일한 클러스터에 분류하는 것이다. 또한 적당한 크기의 클러스터를 생성하도록 한다. 이때 각 어플리케이션은 각 요소들에 대한 정의를 필요로 한다. 방향성 그래프와 같이 표현된 하이퍼텍스트에서는 노드사이의 독립 패스 수에 의해 표시되는 연결강도(strength of connection)를 통해 노드간 관련성을 간접적으로 가시화 한다. 이 경우에도 하나의 링크는 노드간 관련성을 보일 필요가 있다. 만일 도큐먼트의 크기가 변화한다거나, 노드의 링크 위치를 변경했지만 현재 위치를 포인팅하고 있다면 클러스터링은 동일하다. 그러나 문맥내의 변화는 결정적인 요소이기 때문에 명확성을 고려해야 한다. 또한 이러한 변화요소를 분석할 수 있는 방법이 필요하다고 인식하여 분석 방법을 고려하였다.

3.1 다중분리

클러스터링 표준 J에 전적으로 의존하는 접근 방법의 경우에 적용했을 때 가정 만족할 만한 방법이다. 이 경우 다음과 같은 범주를 가지게 된다. L=2라고 가정할 때, 여기에는 몇 가지 차별화된 classification이 존재하고 최소값에 가장 근접한 J를 만들게 된다.

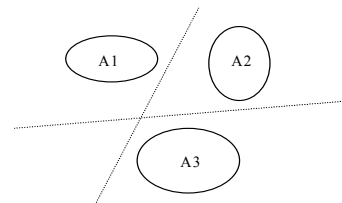


Fig.1 3개 클래스의 2분법

Fig.1은 표본의 모집단을 2개 이분법으로 분리한 것으로 명확하게 3개 클래스를 내포한다. 즉 클래스 L을 포함하는 표본의 모집단이 K 이분법이라고 가정하자. 각 이분법은 2개 그룹으로 표본을 분리한다. 여기에 S_{ij} 는 i번 이분법에 의한 그룹 j에 할당되는 모든 표본의 집합이라고 하고 ($i = 1, 2, \dots, k ; j = 1, 2, \dots, n$) 다음과 같은 2개 조건을 갖는다.

조건1) 하나의 이분은 하나의 그룹 상 각 클래스가 위치한다. 즉 클래스들은 이분에 의해 분리되지 않는다. □

조건2) 클래스 쌍은 적어도 하나의 이분이 있으며, 같은 그룹에 2개의 클래스를 할당하지 않는다. K분

법에 각각에서 하나의 그룹을 선택하고 K 그룹의 교차로서 subset C를 구한다. □

3.3 집합(Aggregation)정의

Def 1

$V' \subset V(G)$ 라 하고 G는 노드의 집합 V' 에 의해 유도된 서브그래프 $\langle V' \rangle$ 는 벡터의 집합 V' 를 갖는다 모든 edge는 V' 에서 양방향 포인트를 갖는다. ■

Def 2

하나의 그래프 G_i 는 최대값으로 서브그래프 G_j 가 없다면 $G_i \subset G_j$ 로 G_i 는 동일 특성을 갖는다. ■

Def 3

하나의 그래프 G의 K-component는 최대값으로 서브그래프 $\langle A \rangle$ 를 유도하며, $\langle A \rangle$ 의 모든 분할 A_1 과 A_2 의 특성을 갖는다. 또한 적어도 $\langle A \rangle$ 의 K개 링크는 A_1 이나 A_2 의 각 객체를 갖는 일련의 이벤트이다. ■

Def 4

집합(Aggregation)을 찾기 위해서 다음과 같은 과정을 거친다. ■

- ① 그룹의 모든 K-component를 정의한다.
- ② 노드를 집합화 한다. 이것은 강 컴포넌트(strong component)를 의미한다.

예) 노드(d,e,f,g,h,i)는 4-component안에 있으며, 다른 강 컴포넌트는 없다.

- ③ 강 컴포넌트는 제거하고, 각 부분들을 모아서 또 다른 약 컴포넌트에 분배한다.

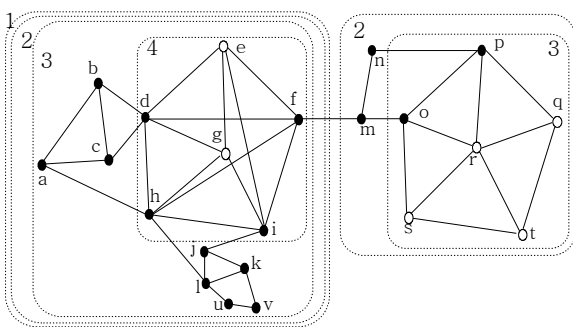


Fig.2 K-component 그래프

예를 들어 2개 노드간 최소 패스 연관이 어떤 것인가를 찾고자 할 경우 b, l노드가 검색되었다고 가정하자, 이때 최소패스는 3으로 b와 i 사이에는 또다른 크기가 3인 패스가 존재한다. 즉 {b,d,g,i}, {b,d,e,i}, {b,d,h,i}, . . . , {d,a,h,i} 사용자는 선택한 패스에 대하여 또다른 경로를 찾을 수 있게된다. 즉 하

나의 그래프는 얻어진 노드사이의 관련성에 관하여 k-component 클러스터링 그래프를 통해서 관련성 강도를 예측할 수가 있다.

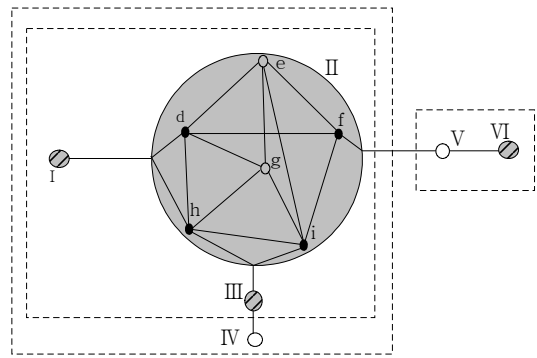


Fig. 3 K-component 클러스터링 그래프

4. 결론

본 논문에서는 하이퍼텍스트로 구성된 검색시스템에 클러스터를 정의하는 방법을 제시하였다. 즉 링크간 독립 패스 수에 있어서 클러스터는 임의적인 특성이 아닌 하이퍼텍스트의 근본적인 특성에 종속적이다. 또한 K-component를 기반으로 하여, 집합화에 의한 클러스터링 방법을 제시하였다. 다만 본 논문에서는 각 집합의 원시의 그래프보다 축소해서 보이므로 복잡하고, 크기가 큰 전체 모습을 보이는데 효과적이다. 또한 집합은 하이퍼텍스트상에서 구조적 질의에 대한 응답에 편리함을 줄수 있으며, 사용자가 여러개 노드를 검색하고자 할 경우, 2개 노드의 관련성을 비교해 볼수 있다.

참고문헌

- [1] F. Harary. Graph Theory. Addison-Wesley, Reading, 1989.
- [2] 이경환, 소프트웨어 재사용을 위한 객체모델링 기법, 교학사, 1993
- [3] Gellersen HW, Gaedke. M , object-oriented Web application development, IEEE Computer society V, v N.1 60-68, 1999.
- [4] D. Budgen and G. Friel, Augmenting the Design Process: Transformations from Abstract Design Representation, CAiSE-92, Springer-verlag, 1992
- [5] Jeffrey M .voas, The challengers of Using COTS SW in Component-based Development, V,31 N 6, 44-45, 1998
- [6] Lisa spicknall Fruth and James M. purtilo, journal of system and software V.32 N.3, 1996.